

Fast Algorithms for the Computation of the Minimum Distance of a Random Linear Code

Fernando Hernando* Francisco D. Igual† Gregorio Quintana-Ortí‡

January 31, 2017

Abstract

The minimum distance of a code is an important concept in information theory. Hence, computing the minimum distance of a code with a minimum computational cost is a crucial process to many problems in this area. In this paper, we present and evaluate a family of algorithms and implementations to compute the minimum distance of a random linear code over \mathbb{F}_2 that are faster than current implementations, both commercial and public domain. In addition to the basic sequential implementations, we present parallel and vectorized implementations that render high performances on modern architectures. The attained performance results show the benefits of the developed optimized algorithms, which obtain remarkable performance improvements compared with state-of-the-art implementations widely used nowadays.

1 Introduction

Coding theory is the area that studies codes with the aim of detecting and correcting errors after sending digital information through an unreliable communication channel. Nowadays, it is widely used in a number of fields, such as data compression [1, 12], cryptography [14, 13], network coding [11], secret sharing [17, 15], etc. The most studied codes are the linear codes, i.e., vector subspaces of dimension k within a vector space of dimension n , and the most used technique to detect and correct errors is via the Hamming minimum distance.

It is well known that if the Hamming minimum distance of a linear code is d , then $d - 1$ errors can be detected, and $\lfloor (d-1)/2 \rfloor$ errors can be corrected. Therefore, it is clear that the knowledge of the minimum distance of a linear code is essential to determine how well such

*Depto. de Matemáticas, Universidad Jaume I, 12.071–Castellón, Spain. carrillf@mat.uji.es.

†Depto. de Arquitectura de Computadores y Automática, Universidad Complutense de Madrid, 28040–Madrid, Spain. figual@uclm.es.

‡Depto. de Ingeniería y Ciencia de Computadores, Universidad Jaume I, 12.071–Castellón, Spain. gquintan@icc.uji.es.

a linear code will perform. This is an active research area because the discovery of better codes with larger distances can improve the recovery and fault-tolerance of transmission lines.

The computation of the minimum distance of a random linear code is an NP-hard problem. Whereas this problem is unsolvable for large dimensions, it can be solved in a finite time for small values of k and n ; in modern computers, feasible values for n are around a few hundreds.

The fastest general algorithm for computing the minimum distance of a random linear code is the so-called Brouwer-Zimmerman algorithm [19], which is described in [5]. This algorithm has been implemented in MAGMA [2] over any finite field, and in GAP (package GUAVA) [3, 4] over fields \mathbb{F}_2 and \mathbb{F}_3 .

Since larger minimum distances allow to detect and recover from larger errors, the interest in computing and designing better linear codes with larger minimum distances is very high. The web page in [6] stores the best minimum distance known to date for every dimension (k and n).

In this paper, we present a family of new efficient algorithms and implementations to compute the minimum distance of random linear codes over \mathbb{F}_2 . Although our implementations only work for \mathbb{F}_2 , our ideas can be applied to other finite fields. The key advantage of the new algorithms is twofold: they perform fewer row additions than the traditional ones, and they increase the ratio between row additions and row accesses. Our new implementations are faster than current ones, both commercial and public-domain, when using only one CPU core. Besides, our new implementations can also take advantage of modern architectures with multiple cores. In this case, our new implementations also render much higher performances than available implementations on multicore architectures. Moreover, our new implementations can take advantage of SIMD (single-instruction, multiple-data) instructions available in old and modern processors to attain even higher performances. With all these improvements, the time required to compute the minimum distance of a linear code can be drastically reduced with respect to current implementations. We believe that the scientific community can benefit from our work because the minimum distance of random linear codes can be computed much faster on serial and multicore architectures by using our new algorithms and implementations.

The rest of the paper is structured as follows. Section 2 introduces the mathematical background of the problem, and classic approaches typically used to solve it (the Brouwer-Zimmerman algorithm), together with the computational challenges related to it. Section 3 proposes a family of new algorithms that tackle the same problem, showing them in an incremental fashion to obtain optimized algorithmic schemes that reduce their computational complexity and thus improve performance on modern computing architectures. Section 4 reports implementation details to adapt the aforementioned new algorithms to modern parallel computing architectures. Section 5 analyzes the performance attained by our new algorithms on different modern architectures, comparing them with state-of-the-art alternative implementations. In Section 6, we leverage our new algorithms to generate some new linear codes. Section ?? briefly describes the programming codes developed in this

paper, their license, and how to get them. Finally, Section 7 closes the paper with some concluding remarks.

2 Background

Let $q = p^r$ be a prime power, we denote by \mathbb{F}_q the finite field with q elements. By a linear code over \mathbb{F}_q , say C , we mean a k -dimensional \mathbb{F}_q vector subspace of \mathbb{F}_q^n , $n > k$, i.e., it is the image of an injective map $i : \mathbb{F}_q^k \hookrightarrow \mathbb{F}_q^n$. As any linear map, i is given by a $k \times n$ matrix G whose rows are a base of C . In coding theory, G is called a generator matrix. Since it has rank k , it can be written, after permutation of columns and elementary row operations, in the systematic form $G = (I_k \mid A)$, where A is a $k \times (n - k)$ matrix, and I_k is the identity matrix of dimension k .

Therefore, if we want to encode a sequence of k bits, say (c_1, \dots, c_k) , we simply multiply $(c_1, \dots, c_k)(I_k \mid A) = (c_1, \dots, c_k, c_{k+1}, \dots, c_n)$. Thus, we introduce $n - k$ redundancy bits, which eventually will be useful to correct the information in case of corruption. A linear code over \mathbb{F}_q contains q^k codewords. When a codeword (c_1, \dots, c_n) is sent through a noisy channel, an error might appear in the received word $r = (c_1, \dots, c_n) + (e_1, \dots, e_n)$, where $e = (e_1, \dots, e_n)$ is the error vector that occurred. The method to find out which codeword was sent when r is received is to replace r with the nearest codeword. In order to do so, we need a metric. Given two vectors in \mathbb{F}_q^n , say $a = (a_1, \dots, a_n)$ and $b = (b_1, \dots, b_n)$, we define the Hamming distance of a and b as the number of positions where they differ, i.e.,

$$d(a, b) = \#\{i \mid a_i \neq b_i\}.$$

Hence, the minimum distance of a linear code is defined as follows:

$$d(C) = \min\{d(a, b) \mid a, b \in C\}.$$

It is quite useful to define the weight of a vector $a = (a_1, \dots, a_n) \in \mathbb{F}_q^n$ to be the number of non-zero positions, i.e.,

$$\text{wt}(a) = \#\{i \mid a_i \neq 0\} = d(a, 0).$$

So, the minimum weight of a linear code is:

$$\text{wt}(C) = \min\{\text{wt}(a) \mid a \in C\}.$$

Since C is a linear subspace, it is easy to prove that $\text{wt}(C) = d(C)$, but computing the weight requires q^k measurements whereas computing the minimum distance requires $\binom{q^k}{2}$. A linear code C over \mathbb{F}_q has parameters $[n, k, d]_q$ if it has length n , dimension k , and minimum distance d .

From now onwards, we will use d instead of $d(C)$ when C is known. With it, we represent either the minimum distance or the minimum weight. This number is essential because $\lfloor (d - 1)/2 \rfloor$ is the number of errors that can be corrected using the nearest codeword

method. If the received word is equidistant to two or more codewords, then we cannot decide which one of them was the sent one. But as far as $\text{wt}(e) \leq \lfloor (d-1)/2 \rfloor$, the nearest codeword to the received one is unique.

Therefore, computing the minimum distance of a linear code is an important task but also a complex one. Actually, Vardy [18] proved that it is an NP-hard problem, and the corresponding decision problem is NP-complete.

The fastest general algorithm for computing the minimum distance of a linear code (to our knowledge) is the Brouwer-Zimmerman algorithm [19], which is explained in detail in [5]. It is implemented, with some improvements, in MAGMA [2] over any finite field. It is also implemented in GAP (package GUAVA) [3] over fields \mathbb{F}_2 and \mathbb{F}_3 .

The method by Brouwer-Zimmerman is outlined in Algorithm 1. It is based on the so called information sets. Given a linear code C with parameters $[n, k, d]$ and a generator matrix G , an information set $S = \{i_1, \dots, i_k\} \subset \{1, \dots, n\}$ is a subset of k indices such that the corresponding columns of G are linearly independent. Therefore, after permutation of columns and elementary row operations we get a systematic matrix $\Gamma_1 = (I_k \mid A_1)$. Assume that we are able to find $m-1$ disjoint information sets ($S_1 \cap \dots \cap S_{m-1} = \emptyset$), then we get $m-1$ different matrices $\Gamma_j = (I_k \mid A_j)$. Notice that there still may be left $n - k(m-1)$ positions, so that the corresponding columns of G do not have rank k but $k_m < k$, then after applying column permutations and row operations, one gets $\Gamma_m = \begin{pmatrix} I_{k_m} & A \\ 0 & B \end{pmatrix}$. In overall, the number of Γ matrices is m : The first $m-1$ will have full rank k , and the last one will have a rank strictly smaller than k .

The idea is to consider an upper bound U , initialized to $n - k + 1$, and a lower bound L , initialized to 1. Then, both bounds are updated after enumerating codewords, and it is checked whether $L \geq U$; if so, the minimum weight is U .

The codewords are enumerated as follows: consider all the linear combinations $c \cdot \Gamma_j$ for $j = 1, \dots, m$, where $c = (c_1, \dots, c_k)$ and $\text{wt}(c) = 1$ (since we are over \mathbb{F}_2 , it means that all c_i are zero but one). After computing any linear combination, if the new weight is smaller than U , then U is updated with the new weight. Moreover, after processing all those linear combinations $c \cdot \Gamma_j$ for $j = 1, \dots, m$, the lower bound is increased in $m-1$ units (actually one after each Γ_j) for the disjoint information sets and a different quantity if the information sets are not disjoint (closed formula). See [5] for more details. Now the same procedure is repeated for linear combinations $c \cdot \Gamma_j$ for $j = 1, \dots, m$ and $\text{wt}(c) = 2$. Then, the same is done for $\text{wt}(c) = 3$, and so on until $L \geq U$ is obtained.

It is clear from line 7 in Algorithm 1 that the lower bound increases linearly in $m-1$ units. So, the more disjoint information sets, the more it increases. Hence, the algorithm will end up earlier with a large m and, in consequence, with a small g . Notice that the number of linear combinations of weight $\leq g$ is $N = \sum_{j=1}^g \binom{k}{j}$. For small values of g we have that $N < \binom{k}{g+1}$, so increasing by one the value of g will require a cost as large as all the previous work. We actually have the following result.

Algorithm 1 MINIMUM WEIGHT ALGORITHM FOR A LINEAR CODE C

Require: The generator matrix G of the linear code C with parameters $[n, k, d]$.

Ensure: The minimum weight of C , i.e., d .

```
1:  $L := 1; U := n - k + 1;$ 
2:  $g := 1;$ 
3: while  $g \leq k$  and  $L < U$  do
4:   for  $j = 1, \dots, m$  do
5:      $U := \min\{U, \min\{\text{wt}(c\Gamma_j) : c \in \mathbb{F}_2^k \mid \text{wt}(c) = g\}\};$ 
6:   end for
7:    $L := (m - 1)(g + 1) + \max\{0, g + 1 - k + k_m\};$ 
8:    $g := g + 1;$ 
9: end while
10: return  $U;$ 
```

Lemma 2.1. *Using the previous notation, if $g \leq \frac{k}{3}$, then*

$$\sum_{j=1}^{g-1} \binom{k}{j} < \binom{k}{g}$$

Proof. First of all we have that $\binom{k}{g} = \frac{k-g+1}{g} \binom{k}{g-1}$. Since $g \leq \frac{k}{3}$ we have that $\binom{k}{g} > 2\binom{k}{g-1} = \binom{k}{g-1} + \binom{k}{g-1} > \binom{k}{g-1} + 2\binom{k}{g-2} > \dots > \sum_{j=1}^{g-1} \binom{k}{j}$. \square

It is quite clear that the most computationally intensive part of this method is the computation of all those linear combinations of the g rows of every matrix Γ_j , since the cost of generating all linear combinations is combinatorial, whereas the cost of the diagonalization to obtain the Γ_j matrices is $O(n^3)$. In the following, we will focus on the efficient computation of the linear combinations to speed up this algorithm.

3 New algorithms

As said before, the most time-consuming part of the Brouwer-Zimmerman algorithm is the generation of linear combinations of the rows of the Γ matrices. Its basic goal is simple: For every Γ matrix, the additions of all the combinations of its rows must be computed, and then the minimum of the weights of those additions must be computed. The combinations of the k rows of the Γ matrices are generated first taking one row at a time, then taking two rows at a time, then taking three rows at a time, etc. The minimum of those weights must be computed, and when the minimum value (L) is equal to or larger than the upper value (U), this iterative process finishes. Unlike the previous algorithm, next algorithms do not show the termination condition and the updatings of L and U in order to simplify the notation.

We have designed several algorithms to perform this task. The basic goals of the new algorithms to obtain better performances are the following:

1. Reduction of the number of row addition operations.
2. Reduction of the number of row access operations.
3. Increase of the ratio between the number of row addition operations and the number of row access operations.
4. Use of cache-friendly data access patterns.
5. Parallelization of the serial codes to use all the cores in the system.
6. Vectorization of the serial and parallel codes to exploit the SIMD/vector hardware machine instructions and units.

In the rest of this section, we describe in detail the algorithms implemented in this work. They are the following: the basic algorithm, the optimized algorithm, the stack-based algorithm, the algorithm with saved additions, and the algorithm with saved additions and unrollings. In those algorithms, we have supposed that vector and matrix indices start with zero.

3.1 Basic algorithm

The most basic algorithm is straightforward: If a Γ matrix has k rows, all the combinations of the k rows taken with an increasing number of rows are generated. For every generated combination, the corresponding rows are added, and the overall minimum weight is updated. The basic algorithm is outlined in Algorithm 2.

Algorithm 2 BASIC ALGORITHM

Require: The generator matrix G of the linear code C with parameters $[n, k, d]$.

Ensure: The minimum weight of C , i.e., d .

```

1: for  $g = 1, 2, \dots$  do
2:   for every  $\Gamma$  matrix ( $k \times n$ ) of  $G$  do
3:     // Process all combinations of the  $k$  rows of  $\Gamma$  taken  $g$  at a time:
4:     ( done, c ) = Get_first_combination();
5:     while ( ! done ) do
6:       Process_combination( c,  $\Gamma$  );
7:       ( done, c ) = Get_next_combination( c );
8:     end while
9:   end for
10: end for
```

Although this algorithm could have been implemented recursively, this type of implementation is usually slow in current machines. Hence, we used instead an iterative implementation.

Each combination c will contain the indices of the rows of the current Γ matrix being processed. The methods `Get_first_combination()` and `Get_next_combination()` return in the first output value whether there are more combinations to process, and compute and return in the second output value the first/next combination to be processed. The second method also requires c as an input parameter to be able to generate the next combination to this one. The order in which the combinations are generated is not important in this algorithm. However, in our implementation the lexicographical order was used since it is very cache-friendly due to accessing the rows in a serial way.

The method `Process_combination` adds the rows of the current Γ matrix with indices in c , and then updates the minimum weight. The algorithm stops as soon as the overall minimum weight is equal or larger than the upper value.

Notice that this algorithm performs $g - 1$ additions of rows for every g rows brought from main memory. Hence, as the ratio row additions/row accesses is so low, it might not be able to extract all the computing power from the processors, and the speed of the main memory could ultimately define the overall performance of the implementation. A first alternative to obtain better performances is to perform fewer additions and accesses. A second alternative is to increase the ratio between row additions and row accesses. Next algorithms will explore both choices.

The basic algorithm is straightforward in its implementation, but it performs many additions of rows. We present the following result.

Lemma 3.1. *The cost in additions of the code inside the innermost **for** loop is:*

$$\binom{k}{g} (g - 1)n$$

Proof. The number of different combinations of k rows taken g rows at a time is $\binom{k}{g}$. For every one of those combinations, the algorithm performs $g - 1$ row additions. Every row addition consists in n additions of bits. By multiplying these three factors, the initial formula is obtained. \square

3.2 Optimized algorithm

In the basic algorithm presented above, the order in which combinations are generated and processed is not very important, except for cache effects. On the other hand, in the optimized algorithm the order in which the combinations are generated is more important, since not all orders can help to reduce the number of additions. The optimized algorithm will use the lexicographical order. In this order, the indices in a combination change from the right-most part. For instance, for 50 elements taken 3 elements at a time, the combinations are generated in the following order: $(0, 1, 2)$, $(0, 1, 3)$, $(0, 1, 4)$, \dots , $(0, 1, 49)$, $(0, 2, 3)$, $(0, 2, 4)$, $(0, 2, 5)$, \dots , $(0, 2, 49)$, $(0, 3, 4)$, \dots

The advantage of the lexicographical order is that each combination is very similar to the previous one. In most cases of this order, there is only one difference between one combination and the next one (or, equivalently, the previous one): the last element. Consequently, in most cases the addition of the first $g-1$ rows performed in one combination can be saved for the computation of the next combination, thus saving many of them. This method is outlined in Algorithm 3.

Algorithm 3 OPTIMIZED ALGORITHM

Require: The generator matrix G of the linear code C with parameters $[n, k, d]$.

Ensure: The minimum weight of C , i.e., d .

```

1: for  $g = 1, 2, \dots$  do
2:   for every  $\Gamma$  matrix ( $k \times n$ ) of  $G$  do
3:     // Process all combinations of the  $k$  rows of  $\Gamma$  taken  $g-1$  at a time:
4:     (done, c) = Get_first_combination();
5:     while ( ! done ) do
6:       Process_all_combinations_starting_with( c,  $\Gamma$  );
7:       (done, c) = Get_next_combination( c );
8:     end while
9:   end for
10: end for

```

The main structure of this algorithm is very similar to the basic one, but there exist two major differences. The first difference is that the combinations are generated with $g-1$ elements instead of g elements. The second difference lies in the processing of the combinations. The method `Process_all_combinations_starting_with` receives a combination c with $g-1$ elements, then it adds the rows with indices in that combination and, finally, it generates all the combinations of g elements that start with the received combination of $g-1$ elements by reusing the previous addition. For instance, if $g = 4$ and $c = (0, 1, 2)$, it will first compute the additions of rows 0, 1, and 2, and then it will reuse that addition to compute the additions of the combinations $(0, 1, 2, 3)$, $(0, 1, 2, 4)$, \dots , $(0, 1, 2, k-1)$, thus saving $g-2$ (2 in this case) additions for every combination.

Lemma 3.2. *The cost in additions of the code inside the innermost **for** loop is:*

$$\left[\sum_{j=g}^{k-1} \left(\binom{j-2}{g-2} (g+k-j-1) \right) \right] n$$

Proof. Assume that c_{g-1} is set to the value j . Then, we have $\binom{j-1}{g-2}$ different combinations on the left part (c_1, \dots, c_{g-2}) . Each of these different combinations require $g-2$ additions and could be combined with any valid value of c_g on the right part, i.e., $k-j$ combinations. So we have $\binom{j-1}{g-2} (g-2+k-j)$ additions. Running through all the possible values of j :

$j = g - 1, \dots, k - 1$, and considering n bit additions per row we have:

$$\left[\sum_{j=g-1}^{k-1} \left(\binom{j-1}{g-2} (g + k - j - 2) \right) \right] n.$$

Modifying the initial index in the summation, we get the initial formula. \square

3.3 Stack-based algorithm

The number of row additions and row accesses can be further reduced by using a stack of $g - 1$ vectors of dimension n . If a combination $c = (c_1, c_2, \dots, c_{g-1})$ is being processed, the stack will contain the following incremental additions: $c_1, c_1 + c_2, c_1 + c_2 + c_3, \dots, c_1 + c_2 + c_3 + \dots + c_{g-1}$. The memory used by the stack is not large at all since both n and g are usually small, and only one bit is needed for each element. In all our experiments, n was always smaller than 300, and the algorithm usually finished for values of g equal to or smaller than 16. (Larger values of g would require a very long time to finish: months or even years of computation in a modern computer.) In overall, in our experiments the stack always used less than 1 KByte of memory.

In the optimized algorithm, the number of additions for every combination of $g - 1$ rows is always the same ($g - 2$ additions) since the addition of the combination is computed from scratch. The desired additions of the combinations of g rows are built on top of those additions with just one extra addition for every combination of g rows.

On the other hand, in the stack-based algorithm, the number of additions for every combination of $g - 1$ rows can be reduced even further if a stack is adequately employed and the combinations are generated in an orderly fashion. In this case the lexicographical order was used, again, because in this order the right-most elements change faster.

As the stack keeps the incremental additions of the previous combination with $g - 1$ elements, it can be used to compute the addition of the current combination with $g - 1$ elements with a lower cost, while at the same time the stack is updated. The number of required additions depends on the left-most element that will change from the previous combination to the current one, since the stack will have to be rebuilt from that level. Hence, to compute a combination of $g - 1$ elements, the minimum number of additions of the new algorithm is one, and the maximum number of additions of the new algorithm is $g - 2$. Consequently, to compute a combination of g elements, the minimum number of additions of the new algorithm is two, and the maximum number of additions of the new algorithm is $g - 1$.

The new stack-based method is outlined in Algorithm 4. Notice how the structure of the new stack-based algorithm is very similar to the previous one. There is a new method called `Initialize_stack` that initializes the necessary data structures to hold the stack. Now, the methods `Get_next_combination` and `Process_all_combinations_starting_with` include a new parameter: the stack with the incremental additions of the rows. The method `Get_next_combination` rebuilds the stack if needed, and therefore it will require the stack as both an input argument and an output argument. It will also require the Γ matrix as an

input argument. On the other hand, the method `Process_all_combinations_starting_with` uses the stack to compute the addition of the rows.

Algorithm 4 STACK-BASED ALGORITHM

Require: The generator matrix G of the linear code C with parameters $[n, k, d]$.

Ensure: The minimum weight of C , i.e., d .

```

1: Initialize_stack( stack );
2: for  $g = 1, 2, \dots$  do
3:   for every  $\Gamma$  matrix ( $k \times n$ ) of  $G$  do
4:     // Process all combinations of the  $k$  rows of  $\Gamma$  taken  $g - 1$  at a time:
5:     ( done, c ) = Get_first_combination();
6:     while ( ! done ) do
7:       Process_all_combinations_starting_with( c, stack,  $\Gamma$  );
8:       ( done, c, stack ) = Get_next_combination( c, stack,  $\Gamma$  );
9:     end while
10:  end for
11: end for

```

Lemma 3.3. *The cost in additions of the code inside the innermost **for** loop is:*

$$\left(\binom{k}{g} + \binom{k-1}{g-1} + \dots + \binom{k-g+2}{2} \right) n$$

Proof. We define the following sets of combinations for $i = 1, \dots, g - 1$:

$$A_i = \{(c_1, \dots, c_g) \mid c_g = c_{g-1} + 1 = c_{g-2} + 2 = \dots = c_{g-i+1} + i - 1\}$$

i.e., at least the last i elements (c_{g-i+1}, \dots, c_g) are consecutive.

Let us consider now a combination (c_1, \dots, c_g) . Assume that the last r elements (c_{g-r+1}, \dots, c_g) are consecutive, but the last $r+1$ are not consecutive, i.e., $c_{g-r} + 1 \neq c_{g-r+1}$. Then, the combination (c_1, \dots, c_g) requires no additions for the rows (c_1, \dots, c_{g-r}) because the contents of the stack is reused. However, r additions are needed to update the higher levels of the stack with (c_{g-r+1}, \dots, c_g) . So, we conclude that a combination with exactly r consecutive elements at the end requires r additions.

The key observation is the following: a combination (c_1, \dots, c_g) with exactly r consecutive elements at the end is contained in exactly r sets A_1, \dots, A_r . Therefore, the number of additions required for a combination is equal to the number of sets A_i where it is contained.

We conclude that the total number of additions required is equal to the sum of the cardinals of the sets A_1, \dots, A_{g-1} . So, in the rest of the proof we are going to calculate the cardinals of these sets.

Notice that A_1 imposes no constraints, i.e., any combination (c_1, \dots, c_g) is contained in A_1 , therefore there are $\binom{k}{g}$ elements in A_1 . Now we consider A_2 . A combination (c_1, \dots, c_g) is in A_2 if and only if $c_g = c_{g-1} + 1$. The latter condition implies that

when (c_1, \dots, c_{g-1}) is fixed, then c_g is automatically fixed. So, we have $\binom{k-1}{g-1}$ different combinations for (c_1, \dots, c_{g-1}) . In general, (c_1, \dots, c_g) is contained in A_i if and only if $c_g = c_{g-1} + 1 = c_{g-2} + 2 = \dots = c_{g-i+1} + i - 1$. The latter condition implies that when (c_1, \dots, c_{g-i}) is fixed, then (c_{g-i+1}, \dots, c_g) is automatically fixed. So, we have $\binom{k-i}{g-i}$ possibilities for (c_1, \dots, c_{g-i}) .

Adding up the cardinals of A_i for $i = 1, \dots, g - 1$ we get the following formula: $\binom{k}{g} + \binom{k-1}{g-1} + \dots + \binom{k-g+2}{2}$. Finally, considering n bit additions per row we get the initial formula of this lemma. \square

3.4 Algorithm with saved additions

The basic algorithm performs $g - 1$ additions to process every generated combination. The optimized algorithm performs only one addition in many cases, and it performs $g - 1$ additions in the rest of cases. The stack-based algorithm performs only one addition in many cases, and it performs between two and $g - 1$ additions in the rest of cases. Now we present an algorithm that performs the same low number of additions ($\lfloor (g - 1)/s \rfloor$) to process every generated combination.

The new algorithm always performs a fixed lower number of additions for all combinations by using a larger additional storage. Its main advantage is its smaller cost. Its main disadvantage is the extra memory space needed, but this issue is not a serious handicap, as the availability of main memory in current computers is usually very large.

For every Γ matrix, this algorithm saves in main memory the additions of all the combinations of the k rows taken g at a time for values of $g = 1, 2, \dots, s$. The value s is fixed at the beginning of the program, and it determines the maximum amount of memory used. In our experiments, we employed values of s up to 5, since it rendered good performances and larger values required too much memory. For instance, with $s = 5$ this algorithm required a storage of around 95 MBytes for processing the new linear codes presented in this paper, which is not an excessive amount, while rendering good performances.

The saved additions of the combinations of k rows taken 1, 2, \dots , s rows at a time will be then used to build the additions of the combinations of k rows taken $s + 1$, $s + 2$, \dots rows at a time. This idea is simple, but the problem is to be able to implement it in a very efficient way.

If these additions are saved in the lexicographical order of the row indices in the combinations, it is really efficient to combine them. This order is key to this algorithm.

Next, we describe some details in the most simple case. If $g = a + b$ with positive numbers a and b such that $a \leq s, b \leq s$, the addition of the rows of the combination c with indices $(c_1, c_2, \dots, c_a, c_{a+1}, \dots, c_g)$ can be computed as the addition of the rows of the following combinations: the combination (c_1, c_2, \dots, c_a) (called left combination) and the combination (c_{a+1}, \dots, c_g) (called right combination). In this way, with just one addition the desired result can be obtained if we have previously saved the additions of the combinations of k rows taken up to at least $\max(a, b)$ at a time.

Therefore, if $g = a + b$, to obtain the combinations of k rows taken g at a time,

the combinations of k rows taken a at a time (left combinations) and the combinations of k rows taken b at a time (right combinations) must be combined. However, not all those combinations must be processed. There is one restriction to be applied to the left combinations, and another one to be applied to the right combinations. Next, both of them are described. Note that it is very important that these restrictions must be applied efficiently to accelerate this algorithm. Otherwise, an important part of the performance gains could be lost.

In the case of the left combinations, not all of the combinations must be processed. For instance, if $k = 50$, $a = 3$, and $b = 2$, left combinations starting with 46 or larger indices should be discarded since no right combination can be appended to form a valid combination (as an example, left combination (46, 47, 48) cannot be concatenated to any right combination of two elements to form a valid combination). If the saved additions of the combinations of k elements taken a at a time are kept in the lexicographical order of the combinations, the following formula returns the index of the first combination in the saved combinations that must not be processed:

$$\binom{k}{a} - \binom{g-1}{a}$$

In the case of the right combinations, not all of the combinations must be processed. The last element in the left combination (c_1, c_2, \dots, c_a) will define the right combinations with which this can be combined, since it can only be combined with combinations starting with $c_a + 1$ or larger values. If the saved additions of the combinations are kept in the lexicographical order of the combinations, we can compute easily which combinations of k rows taken b at a time must be processed. The formula that returns the index of the first right combination to be combined is the following one, where e is the last element in the left combination:

$$\binom{k}{b} - \binom{k-e-1}{b}$$

A general recursive algorithm that works for any k , any g , and any s has been developed. Though recursive algorithms can be slow, ours is really fast because the cost of the tasks performed inside each call is high, and the maximum depth of the recursion is $\lceil g/s \rceil$. The general method is outlined in Algorithm 5.

The data structure that stores the saved additions of the combinations of the rows of every Γ matrix must be built in an efficient way. Otherwise, the algorithm could underperform for matrices that finish after only a few generators. For every Γ matrix, this data structure contains several levels ($l = 1, \dots, s$), where level l contains all the combinations of the k rows of the Γ matrix taken l at a time. The way to do it in an efficient way is to use the previous levels of the data structure to build the current level of the data structure. In our algorithms, to build level l , levels $l-1$ and 1 were used. This combination must be performed in a way that both keeps the lexicographical order in level l and is efficient.

Algorithm 5 ALGORITHM WITH SAVED ADDITIONS

Require: The generator matrix G of the linear code C with parameters $[n, k, d]$.

Ensure: The minimum weight of C , i. e., d .

```
1: Initialize_data_structures_for_storing_additions( SA );
2: for  $g = 1, 2, \dots$  do
3:   for every  $\Gamma_i$  matrix ( $k \times n$ ) of  $G$  do
4:     if  $g \leq s$  then
5:       Generate and save all combinations of  $g$  rows of  $\Gamma_i$  into  $SA_i$  ;
6:     end if
7:     Process_step(  $SA_i, g, \emptyset$  );
8:   end for
9: end for
10: End of Algorithm

11: Method Process_step(  $SA_i, g, c$  ) :
12:  $a := \min(g, s)$ ;
13:  $b := g - a$ ;
14: if  $a < s$  then
15:   Compute the minimum distance of  $SA_i$  adding  $c$  to suited combinations;
16: else
17:   for  $j = \text{index\_of}( k, a, \text{last\_element\_of}( c ) )$  to  $\text{index\_of}( k, a, k - g )$  do
18:      $e = j$ -th combination saved in  $SA_i$ ;
19:     if  $\text{last\_element\_of}(e) + b < k$  then
20:       Process_step(  $SA_i, b, c + e$  );
21:     end if
22:   end for
23: end if
24: End of Method

25: Function index_of(  $p, q, r$  ) :
26: Return  $\binom{p}{q} - \binom{p-r-1}{q}$ 
27: End of Function
```

Lemma 3.4. *The cost in additions of the code inside the loop in line 3 of Algorithm 5 (equivalent part in the previous algorithms) is:*

$$\binom{k}{g} n \lfloor (g-1)/s \rfloor$$

Proof. Obvious, since this algorithm performs only $\lfloor (g-1)/s \rfloor$ additions per combination. \square

Although this algorithm performs much fewer additions than previous ones, it has one drawback: the amount of data being stored and processed is much larger. Even though

the amount of data being stored is not prohibitive (around 50 MB in our experiments), the processing of those data will produce many more cache misses than previous algorithms. Recall that previous algorithms must just process a few Γ matrices of dimension $k \times n$, whereas this algorithm must process a few matrices of dimension $\binom{k}{g} \times n$, for $g = 1, \dots, s$. In the first case, those few Γ matrices can be stored in the first levels of cache memory, whereas in the second case the matrices with the combinations will not usually fit there.

3.5 Algorithm with saved additions and unrollings

All the algorithms described above perform $g - 1$ row additions for every g rows brought from main memory, being the only difference among them the number of total operations. The goal of all of them, except the basic one, is to reduce the total number of row additions and row accesses. Since the ratio row additions/row accesses is so low (close to one), and as main memory is much slower than computing cores, this low ratio might reduce performances when the memory system is specially slow or saturated.

We have developed an algorithmic variant of the algorithm with saved additions that improves performance by increasing the ratio additions/memory accesses, in order to be less memory-bound. The main and only difference is that several combinations are processed at the same time, and whenever one row is brought from main memory, it will be reused as much as possible in order to decrease the number of memory accesses. This technique is called *unrolling*, and it is widely used in high-performance computing. This technique will reduce the number of memory accesses, and consequently the number of cache misses since data are reused when transferred from main memory.

For instance, by processing two combinations at the same time, the number of rows accessed can nearly be halved since each accessed row is used twice, thus doubling the ratio row additions/row accesses. Processing three iterations at a time would improve this ratio even further. If more row additions per every row access are performed, the fast computing cores will work closer to their limits, and main memory will be removed as the limiting performance factor.

In our experiments we have tested the processing of two combinations at a time, and the processing of three combinations at a time. We have not evaluated higher numbers because of the diminishing returns.

The loop in line 17 of algorithm 5 processes one combination in each iteration of the loop. To process two combinations at a time, this loop should be modified to process two iterations of the old loop in each iteration of the new loop.

But executing two or more iterations at a time is more effective when the last element of them are exactly the same. When the last element is the same, both left combinations must be combined with the same subset of right combinations. In contrast, when the last element is different, each left combination must be combined with a different subset, and therefore it is not so effective to blend them. For instance, left combinations $(0, 1, 4)$ and $(0, 2, 4)$ can be executed at the same time, whereas in left combinations $(0, 1, 2)$ and $(0, 1, 3)$ is not so effective.

If the lexicographical order is employed, the last element of the combinations is the one that changes most. In this case, consecutive combinations usually contain different last elements, and the unrolling will not be so effective. Therefore, a new order must be used. The requirements for this new order are two-fold: The first one is that the first element must be the one that changes least so that we can efficiently access all combinations starting with some given element (the last one of the previous combination plus one). The second one is that the last element must change very little, in order to be able to blend as many consecutive combinations as possible. Hence, the order we have used is a variation of the lexicographical one in which the element that changes least is the first one, then the last one, and then the rest. For instance, if $k = 5$ and $g = 3$, the order is the following one: $(0, 1, 2)$, $(0, 1, 3)$, $(0, 2, 3)$, $(0, 1, 4)$, $(0, 2, 4)$, $(0, 3, 4)$, etc.

4 Implementation and optimization details

4.1 Parallelization of the basic algorithm

The outer loop `For g` (line 1 of Algorithm 2) cannot be parallelized since its number of iterations is not known *a priori*. Recall that the iterative process can finish after any iteration of the g loop (whenever the minimum weight is larger than some value). In addition, the cost of every iteration of the g loop is extremely different (the cost of each iteration is usually larger than the cost of all the previous iterations). Furthermore, g is usually a small number. For instance, in our experiments, g was always smaller than or equal to 16. As the number of cores can be higher, the parallelization of this loop would not take advantage of all the computer power. Because of all of these causes, the parallelization of this loop must be discarded.

The middle loop `For every Γ` (line 2 of Algorithm 2) can be easily parallelized (by assigning a different Γ matrix to every core). Despite the cost of processing every Γ matrix is very similar, in order to be able to parallelize this algorithm there should be as many or more Γ matrices than computing cores. However, there are usually many more cores than Γ matrices. For example, in our most time-consuming experiments there were 5 Γ matrices, whereas current computers can have a larger number of cores. So this solution must also be discarded due to its inefficiency and lack of potential scalability.

The `while` loop (line 5 of Algorithm 2) can be easily parallelized, but it has an important drawback that makes the parallelization inefficient: To parallelize that loop, the invocation of `Get_next_combination` (line 7) should be inside a critical region so that only one thread can execute this method at a time. Otherwise, two threads could get the same combination or, worse yet, one combination could be skipped. This parallelization strategy works fine for a very low number of threads (about 2), but when using more threads, the method `Get_next_combination` becomes a big bottleneck, and performances drop significantly.

In conclusion, despite how simple the structure of the basic algorithm is, its parallelization will not usually render high performances. Even in some cases its parallelization could render much lower performances than the serial algorithm.

4.2 Parallelization of the optimized algorithm

As the structure of the optimized algorithm is so similar to the structure of the basic algorithm, its parallelization is going to present the same drawbacks. Although the concurrent method `Process_all_combinations_starting_with` has a larger cost than the analogous one of the previous algorithm, the critical region in method `Get_next_combination` would make performances drop when the number of cores is slightly larger. Therefore, the parallelization of this algorithm must be discarded too.

4.3 Parallelization of the stack-based algorithm

The structure of the stack-based algorithm is very similar to the previous one, and therefore it has the same drawbacks. Its main difference with the optimized algorithm is that the cost of the method `Get_next_combination` is larger since the stack must be rebuilt in some cases. This fact makes the parallelization of this algorithm even less appropriate since the cost of the critical region is larger.

4.4 Parallelization of the algorithms with saved additions

The structures of the algorithm with saved additions and the algorithm with saved additions and unrollings are very similar. So, both of them will be tackled at the same time.

The parallelization of the algorithms with saved additions is very different from the previous ones. As the combinations are already generated and the additions of those are saved in vectors, its parallelization does not require the use of a large critical section, thus rendering higher potential gains in the parallel implementations. The loop that must be parallelized is the `for` loop of the `else` branch, but it must only be parallelized for the first level of the recursion. Hence, we have an algorithm, the algorithm with saved additions, that is both efficient and parallel.

In the parallelization of this algorithm, a dynamic scheduling strategy must be used since the cost of processing every element of the vector (subcombination) is very different. For instance, if $k = 50, g = 6, s = 3$, processing iteration $(0, 1, 2)$ would require much longer than processing iteration $(0, 1, 46)$. In the first case, many combinations must be processed: $(0, 1, 2, 3, 4, 5), (0, 1, 2, 3, 4, 6), \dots$. In the second case, the only combination to be processed would be: $(0, 1, 46, 47, 48, 49)$. In our parallelized codes, we used OpenMP [16] to achieve the dynamic scheduling strategy of mapping loop iterations to cores.

We used a small critical section for updating the overall minimum weight. However, we minimized the impact of this critical section by making every thread work with local variables, and by updating the global variables just once at the end.

4.5 Vectorization and other implementation details

Usual scalar instructions allow to process one byte, one integer, one float, etc. at a time. In contrast, hardware vector instructions allow to process many numbers at a time, thus

improving performances. Vector instructions use vector registers, whose size depends on the architecture. For example, while older x86 architectures used 128-bit registers, modern architectures use 256-bit or even 512-bit registers. With so wide vector registers, just one vector instruction can process a considerable number of elements.

However, one drawback of the vector instructions is its standardization. Every new architecture contributes new vector instructions, and older architectures do not support the new vector instructions. Therefore, developing a vector code is not straightforward because it will depend on the architecture.

One of the commercial and most-employed implementations, MAGMA, only allows the use of hardware vector instructions on the newest architectures with AVX support, and not on processors with SSE. Unlike MAGMA, our implementations are more flexible, and they can use hardware vector instructions both in processors with SSE and AVX support, that is, in both old and new processors, both from Intel and AMD.

When developing a high-performance implementation, the algorithm is very important, but then even some small implementation choices can greatly affect the performances.

We used the C language since it is a compiled language, and therefore it usually renders high performances. We have also chosen it because of its high portability.

In our implementations, we used the 32-bit integer as the basic datatype, thus packing 32 elements into each integer. This provides high performances for the scalar implementations since just one scalar instruction can process up to 32 elements. We also tested 64-bit integers since in one operation more elements would be processed, but performances dropped. We think that the cause of this drop is the additional elements that must be processed when n is not a multiple of the number of bits of the basic datatype. When 64-bit integers are used, $\text{mod}(n/64)$ additional elements must be processed. On the other hand, when 32-bit integers are used, $\text{mod}(n/32)$ additional elements must be processed, which is usually a smaller overhead.

5 Performance analysis

This section describes and analyzes the performance results attained by our implementations, comparing them with state-of-the-art software that perform the same tasks. The experiments reported in this article were performed on the following two computing platforms:

- **Cplex:** This computer was based on AMD processors. It featured an AMD Opteron™ Processor 6128 (2.0 GHz), with 8 cores. Its OS was GNU/Linux (Version 3.13.0-68-generic). Gcc compiler (version 4.8.4) was used. In the experiments we usually used up to 6 cores (of the 8 cores it had) since we were not the only users. Note that some of the experiments presented in this section lasted for weeks.
- **Marbore:** This computer was based on Intel processors. It featured two Intel Xeon® CPUs E5-2695 v3 (2.30 GHz), with 28 cores in total. Its OS was GNU/Linux (Version 2.6.32-504.el6.x86_64). Gcc compiler (version 4.4.7) was used. As this was a dedicated

computer, we used all its cores in the experiments. In this computer the so-called *Turbo Boost* mode of the CPU was turned off in our experiments.

Our implementations were compared with the two most-used implementations currently available:

- MAGMA [2]: MAGMA is a commercial software package designed for computations in algebra, algebraic combinatorics, algebraic geometry, etc. Given the limitations in its license, it was installed only in the AMD computer, and it was not in the Intel computer. MAGMA Version 2.22-3 was employed in our experiments.

As MAGMA only allows the use of hardware vector instructions on the newest processors with AVX support, and not on processors with SSE, we could not use this feature in the AMD computer. Therefore, the version of MAGMA we evaluated only used scalar instructions. On the contrary, our implementations are more flexible, and they can use hardware vector instructions both on processors with SSE and AVX support, that is, on both old and new processors, both from Intel and AMD.

We evaluated both serial and parallel MAGMA since it allows the use of both one and several cores on multi-core architectures.

- GUAVA [3, 4]: GAP (*Groups, Algorithms, Programming*) is a software environment for working on computational discrete algebra and computational group theory. It includes a package named GUAVA that contains software to compute the minimum weight of linear codes. It is public domain and free, and thus it will be evaluated on both architectures. GUAVA Version 3.12 and GAP Version 4.7.8 were employed in our experiments.

Since GUAVA does not allow the use of hardware vector instructions, it only uses scalar instructions. As GUAVA does not allow the use of multiple cores, we evaluated it only on one core. Unlike GUAVA, our implementations can use both scalar and hardware vector instructions, both on one and on several cores.

In the first subsection of this section we will compare the algorithms described in this paper. In the second subsection of this section we will compare the best implementations developed in this paper and the implementations available on one core. In the third subsection of this section we will compare the best implementations developed in this paper and the implementations available on multicore machines. In the fourth section we will explore the scalability and parallel performance of the implementations.

5.1 A comparison of the algorithms and implementations described in this paper

Table 1 reports the time spent by the algorithms to compute the minimum distance of a random linear code with parameters $[150,50,28]$. The performances are very encouraging

since the best algorithm is more than 6 times as fast as the worst one. This improvement reflects the qualities and virtues of some of our algorithms. The optimized algorithm improves the performances of the basic one since it performs fewer additions. The stack-based algorithm improves the performances of the optimized one since it performs even fewer additions. The method based on storing combinations performs even fewer additions by storing and reusing previous additions, and thus it renders higher performances. The algorithm with saved additions and unrollings performs exactly the same additions as the algorithm with saved additions, but it performs fewer memory accesses, thus attaining better performances in one of the two computers: the one with the slowest memory system. For larger codes the performances of the last algorithm were better than those shown in the table, since the number of saved data was larger and the memory started to become a serious bottleneck.

Implementation	cplex	marbore
Basic	511.5	273.1
Optimized	191.0	105.3
Stack-based	144.5	84.6
Saved additions with $s = 5$	89.8	44.7
Saved additions with $s = 5$ and unrollings by 2	74.0	44.6

Table 1: Time (in seconds) for the different implementations on a linear code with parameters $[150, 50, 28]$.

5.2 A comparison of the best implementations on one core

Table 2 compares our best implementations and the two best implementations available, MAGMA and GUAVA, for some linear codes of medium size on one core of `cplex`.

As MAGMA cannot use hardware vector instructions in the computer used in the experiments, we could only evaluate it with scalar instructions. As GUAVA cannot use hardware vector instructions at all, we could only evaluate it with scalar instructions.

We show two implementations of our best algorithm: the first one is the usual implementation that only uses scalar instructions, and the second one is an implementation that uses hardware vector instructions.

Both of our new implementations clearly outperform the other two in all cases. Our new implementations are several times faster than the current ones. The performance improvement is remarkable in these cases: Our scalar implementation is in average 3.24 times as fast as MAGMA, and our vector implementation is in average 5.75 times as fast as MAGMA. Our scalar implementation is in average 2.62 times as fast as GUAVA, and our vector implementation is in average 4.67 times as fast as GUAVA.

Table 3 compares the implementations for the new linear codes on one core of `cplex`. These linear codes are larger than the previous linear ones. Our two new implementations

Code	MAGMA	GUAVA	Scalar Saved	Vector Saved
[150, 50, 28]	161.1	193.6	74.0	38.9
[130, 67, 15]	1,980.0	1,585.5	574.8	331.3
[115, 63, 11]	5,056.7	3,703.7	1,292.8	755.1
[102, 62, 12]	20,585.9	14,356.8	5,258.0	3,047.5
[150, 77, 17]	53,052.9	40,804.3	19,262.2	10,245.4

Table 2: Time (in seconds) for the best implementations on several linear codes of medium size on one core of `cplex`.

clearly outperform the usual current implementations. Our scalar implementation is in average 1.71 times as fast as MAGMA, and our vector implementation is in average 3.58 times as fast as MAGMA. Our scalar implementation is in average 1.64 times as fast as GUAVA, and our vector implementation is in average 3.44 times as fast as GUAVA.

Code	MAGMA	GUAVA	Scalar Saved	Vector Saved
[235, 51, 64]	802,364.2	681,339.1	484,788.1	234,890.8
[236, 51, 64]	786,181.6	686,686.4	484,834.9	231,345.5
[233, 51, 62]	643,663.4	567,629.1	335,678.6	159,470.8
[233, 52, 61]	687,073.4	934,105.0	482,535.2	225,141.2
[232, 51, 61]	503,984.2	456,413.2	261,250.3	125,695.6

Table 3: Time (in seconds) for the best implementations on the new linear codes on one core of `cplex`.

Figure 1 shows the performances (in terms of combinations per second) for all the linear codes on one core of `cplex`. Therefore, the higher the bars, the better the performances are. The total number of combinations processed by all the algorithms are usually similar, but not identical, since some algorithms (such as GUAVA) can process an additional Γ matrix in a few cases. To make the comparison fair, all the algorithms used the same total number of combinations: those returned by GUAVA. This figure shows that both the new scalar and vector implementations clearly outperform both MAGMA and GUAVA for all cases.

5.3 A comparison of the best implementations on multicores

Table 4 compares our best implementations and the best implementations available for some linear codes of medium size on multiple cores of `cplex`.

As GUAVA cannot use several cores, it was not included in these experiments. As MAGMA cannot use hardware vector instructions in the computer being used, we could

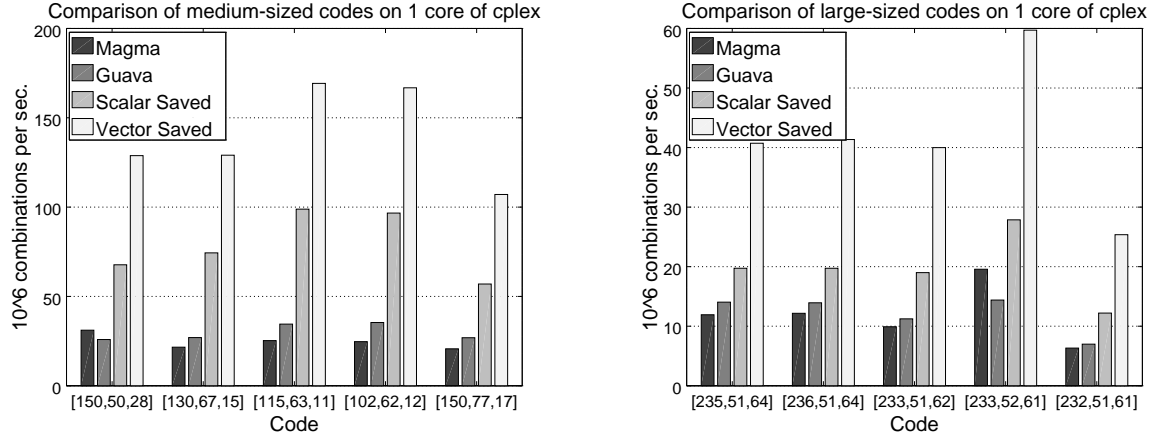


Figure 1: Performance (in terms of 10^6 combinations per sec.) of the best implementations for all the linear codes on one core of `cplex`.

only evaluate it with scalar instructions.

The performance improvement is also remarkable in these cases: Our new implementations are several times faster than MAGMA. Our scalar implementation is in average 3.26 times as fast as MAGMA, and our vector implementation is in average 5.58 times as fast as MAGMA. The improvement factors of our new implementations with respect to MAGMA on multicores are similar to those on a single core, thus showing that our parallelization is as good as that of MAGMA.

Code	MAGMA	Scalar Saved	Vector Saved
[150, 50, 28]	29.3	13.5	7.8
[130, 67, 15]	345.3	104.7	64.9
[115, 63, 11]	860.5	224.4	134.2
[102, 62, 12]	3,659.8	894.7	525.1
[150, 77, 17]	9,562.8	3,314.6	1,763.1

Table 4: Time (in seconds) for the best implementations on several linear codes of medium size on 6 cores of `cplex`.

Table 5 compares our best implementations and the best implementations available for the new linear codes on multiple cores of `cplex`. The performance improvement is also remarkable in these cases: Our new implementations are faster than MAGMA. In average, our scalar implementation is 1.72 times as fast as MAGMA, and our vector implementation is 3.63 times as fast as MAGMA. The improvement factors of our new implementations with respect to MAGMA on multicore are also similar to those on a single core, thus showing that our parallelization is as good as that of MAGMA.

Code	MAGMA	Scalar Saved	Vector Saved
[235, 51, 64]	133,233.3	80,963.0	38,428.5
[236, 51, 64]	132,385.7	80,966.6	38,497.1
[233, 51, 62]	108,552.2	56,083.6	26,725.6
[233, 52, 61]	116,002.6	80,583.9	37,660.1
[232, 51, 61]	85,341.1	43,618.6	20,691.6

Table 5: Time (in seconds) for the best implementations on the new linear codes on 6 cores of `cplex`.

Figure 2 shows the performances (in terms of combinations per second) for all the linear codes on 6 cores of `cplex`. Therefore, the higher the bars, the better the performances are. To make the comparison fair, all the algorithms used the same total number of combinations: those returned by `GUAVA`. This figure shows that both the new multicore scalar and vector implementations clearly outperform `MAGMA` for all cases.

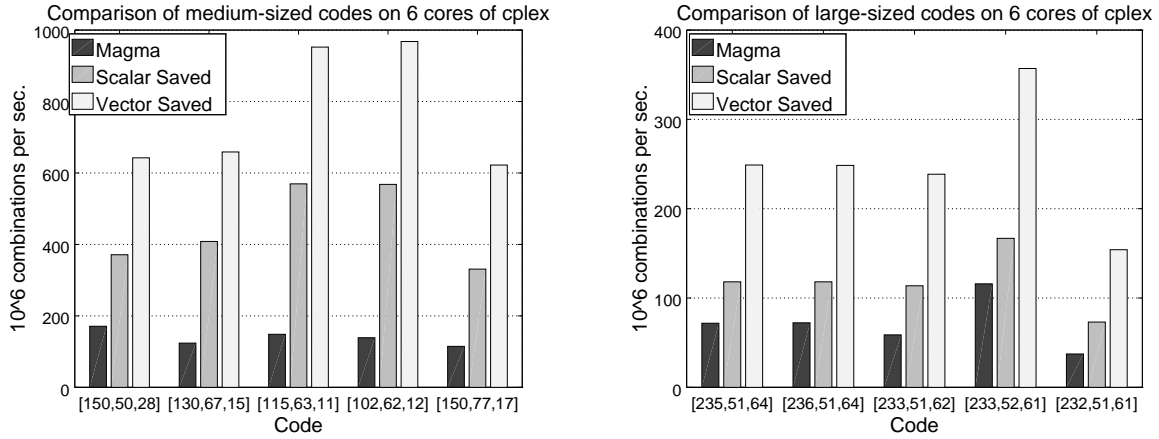


Figure 2: Performance (in terms of 10^6 combinations per sec.) of the best implementations for all the linear codes on 6 cores of `cplex`.

5.4 Parallelization and scalability

Figure 3 shows the obtained speedups by both our scalar and vector algorithms with saved additions for several configurations of cores on both `cplex` and `marbore`. We employed a medium-sized code with parameters [150,77,17], and similar results were obtained on other codes. In this cases, we used up to all the 8 cores in `cplex`. The two plots show that the new implementations are remarkably scalable, even with a high number of cores, since the obtained speedups are very close to the perfect ones. For instance, when run on the 28

cores of **marbore**, the parallel implementation was more than 26 times as fast as the serial implementation.

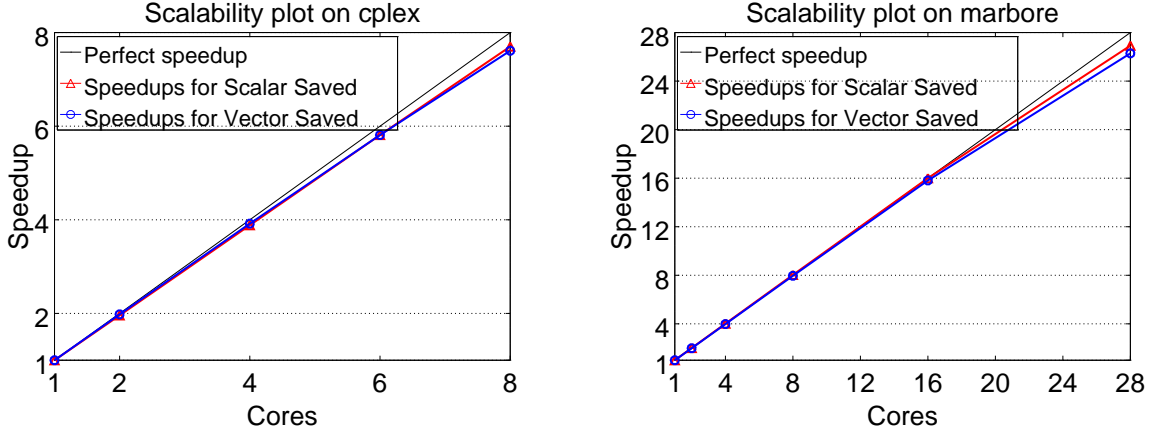


Figure 3: Obtained speedups on both machines (**cplex** left; **marbore** right).

6 New linear codes

Armed with our new implementations, computations that used to take several days using state-of-the-art software packages can be performed in only a few hours. Taking advantage of this fact, we were able to explore well known techniques to generate new linear codes.

We consider matrix-matrix product codes with polynomial units (see [9]) $C = [C_1 C_2] \cdot A$, where C_1 and C_2 are cyclic nested codes with the same length m and $d_2 > 2d_1$, and matrix A is defined as:

$$A = \begin{pmatrix} 1 & p \\ 0 & 1 \end{pmatrix},$$

where p is a unit in the following ring $\mathbb{F}_2[x]/(x^m - 1)$.

We compared the minimum distance of these binary linear codes obtained with our implementations with the ones in [6], the well-known archive of best linear codes. Later on, we obtained the following linear codes whose parameters are better than the ones previously known in [6]:

From [6]	New codes
[234, 51, 62]	$C_1 = [234, 51, 63]$
[234, 52, 61]	$C_2 = [234, 52, 62]$

$C_1 = [C_1, C_2] \cdot A$, where $C_1 = (f_1)$ and $C_2 = (f_2)$ with:

- $f_1 = x^{67} + x^{59} + x^{54} + x^{51} + x^{49} + x^{42} + x^{39} + x^{36} + x^{35} + x^{34} + x^{33} + x^{31} + x^{30} + x^{29} + x^{27} + x^{26} + x^{25} + x^{24} + x^{22} + x^{21} + x^{19} + x^{17} + x^{16} + x^{15} + x^{14} + x^{13} + x^{11} + x^6 + x^5 + x^3 + x^2 + 1,$
- $f_2 = (x^{117} - 1)/(x + 1),$
- $p = x^{117} + x^{116} + x^{115} + x^{111} + x^{110} + x^{109} + x^{103} + x^{102} + x^{98} + x^{95} + x^{94} + x^{92} + x^{88} + x^{85} + x^{83} + x^{81} + x^{74} + x^{72} + x^{70} + x^{68} + x^{66} + x^{65} + x^{64} + x^{62} + x^{61} + x^{58} + x^{56} + x^{55} + x^{54} + x^{51} + x^{49} + x^{45} + x^{43} + x^{39} + x^{38} + x^{37} + x^{36} + x^{35} + x^{34} + x^{28} + x^{27} + x^{23} + x^{22} + x^{20} + x^{19} + x^{16} + x^{14} + x^9 + x^7 + x^6 + x^5 + x^4 + x^3 + x^2 + 1.$

$\mathcal{C}_2 = [C_1, C_2] \cdot A$, where $C_1 = (f_1)$ and $C_2 = (f_2)$ with:

- $f_1 = x^{67} + x^{59} + x^{54} + x^{51} + x^{49} + x^{42} + x^{39} + x^{36} + x^{35} + x^{34} + x^{33} + x^{31} + x^{30} + x^{29} + x^{27} + x^{26} + x^{25} + x^{24} + x^{22} + x^{21} + x^{19} + x^{17} + x^{16} + x^{15} + x^{14} + x^{13} + x^{11} + x^6 + x^5 + x^3 + x^2 + 1,$
- $f_2 = (x^{117} - 1)/(x^2 + x + 1),$
- $p = x^{217} + x^{214} + x^{213} + x^{211} + x^{210} + x^{209} + x^{207} + x^{205} + x^{203} + x^{202} + x^{200} + x^{198} + x^{195} + x^{194} + x^{193} + x^{192} + x^{190} + x^{189} + x^{186} + x^{185} + x^{183} + x^{182} + x^{180} + x^{176} + x^{175} + x^{173} + x^{172} + x^{171} + x^{169} + x^{167} + x^{165} + x^{164} + x^{161} + x^{160} + x^{159} + x^{155} + x^{154} + x^{151} + x^{150} + x^{148} + x^{147} + x^{144} + x^{143} + x^{142} + x^{141} + x^{140} + x^{137} + x^{135} + x^{132} + x^{130} + x^{129} + x^{128} + x^{127} + x^{125} + x^{124} + x^{122} + x^{121} + x^{119} + x^{118} + x^{116} + x^{112} + x^{107} + x^{105} + x^{103} + x^{102} + x^{99} + x^{97} + x^{90} + x^{89} + x^{88} + x^{87} + x^{82} + x^{76} + x^{74} + x^{71} + x^{69} + x^{68} + x^{66} + x^{64} + x^{60} + x^{53} + x^{51} + x^{50} + x^{47} + x^{45} + x^{43} + x^{40} + x^{39} + x^{37} + x^{36} + x^{35} + x^{34} + x^{33} + x^{31} + x^{30} + x^{29} + x^{28} + x^{26} + x^{24} + x^{21} + x^{20} + x^{18} + x^{17} + x^{15} + x^{14} + x^{12} + x^8 + x^5 + x^4 + x^3 + 1.$

Moreover, operating on \mathcal{C}_1 and \mathcal{C}_2 we got five more codes:

From [6]	New codes	Method
[235, 51, 62]	$\mathcal{C}_3 = [235, 51, 64]$	Extend Code(\mathcal{C}_1)
[236, 51, 63]	$\mathcal{C}_4 = [236, 51, 64]$	Extend Code(\mathcal{C}_3)
[233, 51, 61]	$\mathcal{C}_5 = [233, 51, 62]$	Puncture Code($\mathcal{C}_1, 234$)
[232, 51, 60]	$\mathcal{C}_5 = [233, 51, 61]$	Puncture Code($\mathcal{C}_1, 234, 233$)
[233, 52, 60]	$\mathcal{C}_5 = [233, 52, 61]$	Puncture Code($\mathcal{C}_2, 234$)

7 Conclusions

In this paper, we have presented several new algorithms and implementations that compute the minimum distance of a random linear code over \mathbb{F}_2 . We have compared them with the existing ones in MAGMA and GUAVA in terms of performance, obtaining faster implementations in both cases using both sequential and parallel implementations, each of them either in the scalar or in the vectorized case. Finally, we have used our implementation to find out new linear codes over \mathbb{F}_2 with better parameters than the currently existing ones.

The new ideas and algorithms introduced in this paper can also be extended and applied over other finite fields.

Future work in this area will investigate the development of specific new algorithms and implementations for new architectures such as distributed-memory architectures and GPGPUs (General-Purpose Graphic Processing Units).

Source code availability

The source codes described in this paper can be obtained by sending an email to `gquintan@icc.uji.es`, and will be of public access upon the acceptance of this paper. We provide the source codes as we think the scientific community can benefit from our work by being able to compute minimum distances of random linear codes in a faster way.

Acknowledgements

This work is supported by the Spanish Ministry of Economy (grants MTM2012-36917-C03-03 and MTM2015-65764-C3-2-P) and by the University Jaume I (grant P11B2015-02 and TIN2012-32180).

The authors would like to thank Claude Shannon Institute for granting access to `Cplex`.

References

- [1] Ancheta, T. *Syndrome-source-coding and its universal generalization*. IEEE Transactions on Information Theory, 1976, 22, 4, pp. 432–436.
- [2] Bosma, W., Cannon, J., Playoust, C. *The Magma algebra system. I. The user language. Computational algebra and number theory*. (London, 1993). J. Symbolic Comput. 24 (1997), no. 3–4, pp. 235–265.
- [3] The GAP Group. *GAP – Groups, Algorithms, and Programming, Version 4.7.8*. 2015, <http://www.gap-system.org>.
- [4] Baart, R., Boothby, T., Cramwinckel, J., Fields, J., Joyner, D., Miller, R., Minkes, E., Roijackers, E., Ruscio, L. and Tjhai, C. *GUAVA, a GAP package for computing with error-correcting codes, Version 3.12*. 2012, (Refereed GAP package), <http://www.southernct.edu/~fields/Guava/>.
- [5] Grassl, M. *Searching for linear codes with large minimum distance. Discovering mathematics with Magma*. pp. 287–313, Algorithms Comput. Math., 19, Springer, Berlin, 2006.
- [6] Grassl, M. *Bounds on the minimum distance of linear codes*. Online available at <http://www.codetables.de>, 2007. Accessed on 2009-03-27.

- [7] Hernando, F., Hholdt, T., Ruano, D. *List decoding of matrix-product codes from nested codes: an application to quasi-cyclic codes*. Adv. Math. Commun, 6 (2012), pp. 259–272.
- [8] Hernando, F., Lally, K., Ruano, D. *Construction and decoding of matrix-product codes from nested codes*. Appl. Algebra Eng. Comm. Comp., 20 (2009), pp. 497–507.
- [9] Hernando, F., Ruano, D. *New linear codes from matrix-product codes with polynomial units*. Adv. Math. Commun., 4 (2010), pp. 363–367.
- [10] Hernando, F., Ruano, D. *Decoding of matrix-product codes*. J. Algebra Appl., 12 (2013), article id. 1250185.
- [11] Li S. Y. R., Yeung R. W., Cai N. *Linear network coding*. IEEE Transactions on Information Theory, 2003, 49, 2, pp. 371–381.
- [12] Liveris A. D., Xiong Z., Georghiades C.N. *Compression of binary sources with side information at the decoder using LDPC codes*. IEEE Communications Letters, 2002, 6, 10, pp 440–442.
- [13] McEliece, R. J. *A Public-Key Cryptosystem Based On Algebraic Coding Theory*. Deep Space Network Progress Report, 1978, 44, pp. 114–116.
- [14] Niederreiter H. *Knapsack-type cryptosystems and algebraic coding theory*. Problems Control Inform. Theory/Problemy Upravlen. Teor. Inform., 15, 1986, 2, pp. 159–166.
- [15] Olav G., Stefano M., Matsumoto, Ryutaroh M., Diego R., Yuan, L. *Relative generalized Hamming weights of one-point algebraic geometric codes*. IEEE Trans. Inform. Theory, 60, 2014, 10, pp. 5938–5949.
- [16] OpenMP Architecture Review Board. *OpenMP Application Program Interface Version 3.0*. May 2008, <http://www.openmp.org/mp-documents/spec30.pdf>
- [17] Shamir A. *How to share a secret*. Communications of the ACM, (1979) 22, 11, pp. 612–613.
- [18] Vardy, A. *The intractability of computing the minimum distance of a code*. IEEE Trans. Inform. Theory 43 (1997), no. 6, pp. 1757–1766.
- [19] Zimmermann K.-H., *Integral Hecke Modules, Integral Generalized Reed-Muller Codes, and Linear Codes*. Technische Universitat Hamburg-Harburg, Tech. Rep. 3-96, 1996.